

How to Evaluate Your Dialogue Models: A Review of Approaches

Xinmeng Li, Wansen Wu, Long Qin, and Quanjun Yin

College of Systems Engineering, National University of Defense Technology,
Changsha 410000, China
xml.nudt@gmail.com

Abstract. Evaluating the quality of a dialogue system is an understudied problem. The recent evolution of evaluation method motivated this survey, in which an explicit and comprehensive analysis of the existing methods is sought. We are first to divide the evaluation methods into three classes, i.e., automatic evaluation, human-involved evaluation and user simulator based evaluation. Then, each class is covered with main features and the related evaluation metrics. The existence of benchmarks, suitable for the evaluation of dialogue techniques are also discussed in detail. Finally, some open issues are pointed out to bring the evaluation method into a new frontier.

Keywords: Dialogue System · User Simulator · Evaluation Method.

1 Introduction

Automatic evaluation is a non-trivial and challenging task in many sub-fields of natural language processing, e.g. text summarization [33], machine translation [44,34] and dialogue system [37]. Especially, the evaluation of dialogue system often poorly correlates with human judgement. This mismatch is a key bottleneck in migrating dialogue systems developed off-line for application in the real world [37].

The goal of an evaluation method is to assess the performance of a system, which can be defined as “the ability of a system to provide the function it has been designed for” [18]. Prior work often draws on the experience of evaluation methods in other natural language generation tasks, such as BLEU [44], NIST [34] and ROUGE-L [33], which are widely used in machine translation, text summarization and image description. Take advantage of the progress of deep learning, efforts have done to learn evaluation metrics by neural networks [37,60,39,76,75].

In this article, we provide a thorough review of currently available dialogue evaluation paradigms. A method-based taxonomy of the existing evaluation metric is attempted here, including automatic corpus-based evaluation, human-involved evaluation and user simulator based evaluation. We give a condensed overview on Fig.1, where an illustrative representation of the taxonomy is depicted. Deriu et al.[5] has made a survey of evaluation methods for dialogue

systems. In contrast to their work, we focus specifically on generative dialogue evaluation, and make a systematic review of existing techniques, including not only the state-of-the-arts but also those with latest trends. We further introduce how to take advantage of existing benchmarks.

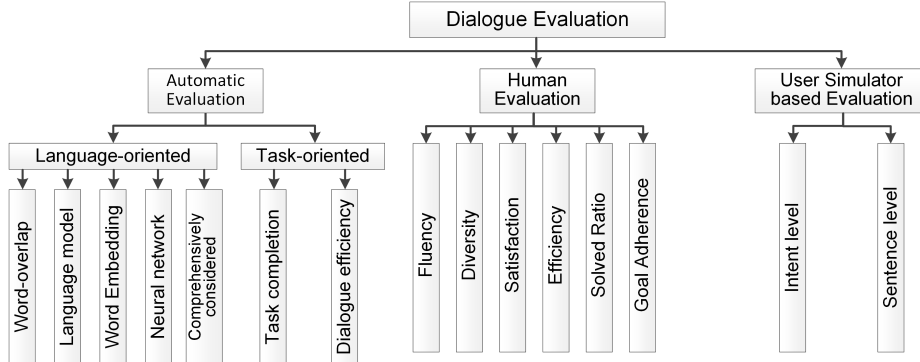


Fig. 1. Taxonomy of the evaluation methods.

2 Overview of Automatic Evaluation Methods

In this section, we introduce the automated approaches for evaluating dialogue systems by dividing them from two aspects, i.e., the language quality and task completion. The language-oriented evaluation reflects the informativity, coherence, fluency and grammaticality of the generated response. From another point of view, the task-oriented evaluation quantify the success rate and dialogue efficiency of a dialogue system.

2.1 Language-oriented evaluation

Word-overlap based metric If there is a standard response in corpus, the simple and natural way is to compare the word-overlap rate of the generated response with the gold one. BLEU [44], NIST [34] and METEOR [1] are initially developed for machine translation evaluation. BLEU computes geometric mean of the precision for n-gram, while NIST replaced geometric mean with arithmetic mean. METEOR considers precision and recall and identifies synonyms and paraphrases between the system output and the ground truth for more comprehensive matching. ROUGE-L [33] is a n-gram based recall which is first used in text summarization task. CIDEr [63] synthetically considers TF-IDF weighting and n-grams averaging. All of them are frequently used word-overlap based metrics and are widely applied to dialogue system and natural language generation tasks [42,69]. Word-overlap based metrics correlate well with human judgements in lower diversity such as machine translation, but often fail in capturing the semantic similarity between the generated sentences and ground truths in dialogue systems for its diversity and dynamic nature. Finch et al.[8] argues that

for most word-overlap metrics, a minimum of 4 references are needed in order to achieve reliable results. These metrics correlate better with human evaluation for datasets providing multiple ground truth sentences, which would lead to bias for reference sentence number [55,15].

Several variants of BLEU are also proposed to improve the evaluation correlation. Sun et al.[58] proposes iBLEU, a revised BLEU score to avoid trivial self-paraphrase and to measure the adequacy and diversity of the generated paraphrase sentence from SMT systems. Galley et al.[11] proposes Δ BLEU to evaluate conversational response generation task that admits a wide variety of possible outputs. However, upfront cost is needed to pay for human rating of the reference set. To remove the human intervention, v BLEU [74] uses a neural network in place of human to annotate the reference set.

Language model based metric Perplexity is widely used for measuring the performance of language models and also applied to conversational models [65,56]. Perplexity can give lower scores to similar dialogues, however the perplexity figure can be difficult to interpret as it ranges from 1 to infinity [46]. The lexical diversity calculates the ratio of unique n-grams to total number of tokens in the dataset which expressing the variety of surface forms as opposed to repeating the same words and phrases [54]. Distinct-1 and distinct-2 are applied to computed the number of distinct uni-grams and bi-grams divided by total number of generated words [27]. Vocabulary size and average utterance length can also be used to measure the language diversity [56].

Embedding based metric Embedding based metrics compute the similarity between of the predicted and the reference sentences by word embeddings. Greedy Matching [47] greedily matches each word in the candidate sentence to a word in the reference sentence based on the cosine similarity of their embeddings. Embedding average [70] computes sentence-level embeddings by averaging the vector representations of their constituent words while Vector Extrema [10] takes the most extreme value of the vector representations of their constituent words for each dimension of the embedding. Different from aforementioned methods obtained through distributional word embeddings, skip-thought vector [23] is base on distributed sentence representations. These metrics are appropriate to evaluate how semantically relevant and on-topic the responses are in evaluating the task of dialogue response generation [53]. However, embedding-based metrics only consist of basic averages of vectors obtained through distributional semantics, they are insufficiently complex for modeling sentence-level compositionality in dialogues [35]. To alleviate the problem, MoverScore [77] and SMS [2] further use sentence and word mover similarity for multi-sentence evaluating.

Neural network based metric Emerging trends are devoted to learning evaluation metrics by neural networks to get rid of the heuristic strategies [37,60,75,52]. Ryan et al.[37] presents ADEM to predict human-like scores for dialogue responses. ADEM computes the score using a dot-product between the vector representations of the model response, the context and the reference response in a linearly transformed space to capture semantic similarity.

BERTScore [75] and BLEURT [52] take advantage of pre-trained language embeddings and computed the similarity score of two sentences as a sum of cosine similarities between their tokens’ embeddings.

Reference-free evaluation methods are also suggested to get rid of ground truth sentences [6]. An adversarial loss could be a way to directly evaluate the extent to which generated dialogue responses sound like they come from a human [21]. RUBER [60], USR [39] and RoBERTa-eval [76] considers both the referenced and unreferenced metrics with heuristic strategies to evaluate the response quality. USR [39], an unsupervised and reference-free evaluation metric is proposed to address the shortcomings of standard metrics for language generation and shown strong correlation with human judgement. RoBERTa-eval [76] investigates to use reference-free metrics, semi-supervised training, and pre-trained text encoders to reduce the bias with human judgement.

Comprehensively considered metric In addition to semantically coincidence, some other metrics such as fluency, topicality and grammaticality also are key factors and should be given attention in evaluation. Readability and grammaticality are considered in [41] for sentence-level NLG evaluation. The evaluation is carried out on Flesch Reading Ease score [9] and Stanford parser score¹. The results illustrate that word-based metrics show better correlations to human ratings of informativeness, whereas grammar-based metrics show better correlations to quality and naturalness. Dziri et al.[7] also includes entailment as an option to approximate dialogue coherence and quality. Guo et al.[14] proposes to evaluate dialog quality with a series of topic-based metrics such as average topic depth, coarse topic breadth, topic keyword frequency and topic keyword coverage to evaluate the ability of conversational bots to lead and sustain engaging conversations on a topic and the diversity of topics the bot can handle. The results show that a user’s satisfaction correlated well with long and coherent on-topic conversations.

2.2 Task-oriented evaluation

Traditional task-oriented dialogue systems often adopt pipeline architecture [73]. There are metrics to evaluate performance of individual modules, such as intent accuracy for the natural language understanding, slot accuracy for dialogue state tracking, task success for policy learning and BLEU for natural language generation [59]. In addition, a comprehensive metric is needed to evaluate holistic system performance. Generally, the task-oriented evaluation focuses on the task completion rate and dialogue efficiency of a dialogue system [66].

Task completion The task-oriented dialogue systems are developed to assist users in achieving specific goals. So the task completion rate, i.e., how well the dialogue system fulfills user’s requirement, is of top priority to evaluate a dialogue system. Entity match rate evaluates task completion by determining if the system generates all correct constraints to search the indicated entities of the user [68]. Task success rate evaluates if the system answered all the associated information

¹ <http://nlp.stanford.edu/software/parser-faq.shtml>

(e.g. phone number, ticket price) [68,26,56,38]. Some fine distinctions of dialogue strategies can also be measured by inappropriate utterance ratio, turn correction ratio, concept accuracy, implicit recovery and transaction success [19,4].

Dialogue efficiency Dialogue efficiency measures the cost incurred in a dialogue, such as the dialogue length or the elapsed time. A direct yet effective metric is the turn number.

Reinforcement Learning based dialogue systems are mainly aimed to optimized task success rate and turn number [31,29,71], where the reward is shaped to better correlate with user satisfaction or be consistent with expert demonstrations directly [30]. Success rate may only measure one aspect of the dialogue policy’s quality. Focusing on information-seeking tasks, Ultes et al.[61] proposes an interaction quality based reward to balance the dialogue policy.

Furthermore, success rate and turn number are also commonly used as system-level evaluation metrics for dialogue policy management and dialogue system performance [57,45].

3 Human evaluation

Automated evaluation methods focus on single-turn quality, so they are often incapable to evaluate the holistic performance of dialogue system. Researchers have found that the actual system-level performances do not match the automatic evaluation results when interacting with real users [3]. For practical use, a well-performed dialogue system is one which can interact efficiently and naturally with human subjects to complete an application-specific task [37]. So an ideal evaluation is to recruit human beings to interact with a dialogue agent to check whether it can successfully assist users to accomplish a given task. Large-scale human evaluation is always performed in crowd-sourcing platform (i.e., CrowdFlower², Amazon Mechanical Turk³), where subjects interact with the dialogue systems and rate them via questionnaires [42,53,54,28].

Human evaluation is a high subjective endeavor, so the measures often differentiate between research groups [20]. Anu et al.[64] proposes a hybrid metric based on engagement, domain coverage, coherence, topical diversity and conversational depth. See et al.[51] designs two controllable neural text generation methods and conducted human evaluation to measure the effect of these control parameters from eight aspects: avoiding repetition, interestingness, making sense, fluency, listening, inquisitiveness, humanness and engagingness. Ghandeharioun et al.[13] uses a comprehensive evaluation strategy to approximate sentiment, semantic similarity, and engagement. Users are asked to give a rating on the naturalness, coherence, and task-completion capability of a system in [32]. Shi et al.[56] asks evaluators to interact with trained systems and obtained their opinions on satisfaction, efficiency, naturalness and rule-likeness and solved Ratio. Experimental results show that the auto-success is not necessarily correlated with the user-rated solved ratio. So human evaluation is indispensable to correct the automatic evaluation bias.

² faircrowd.work/platform/crowdfower

³ www.mturk.com

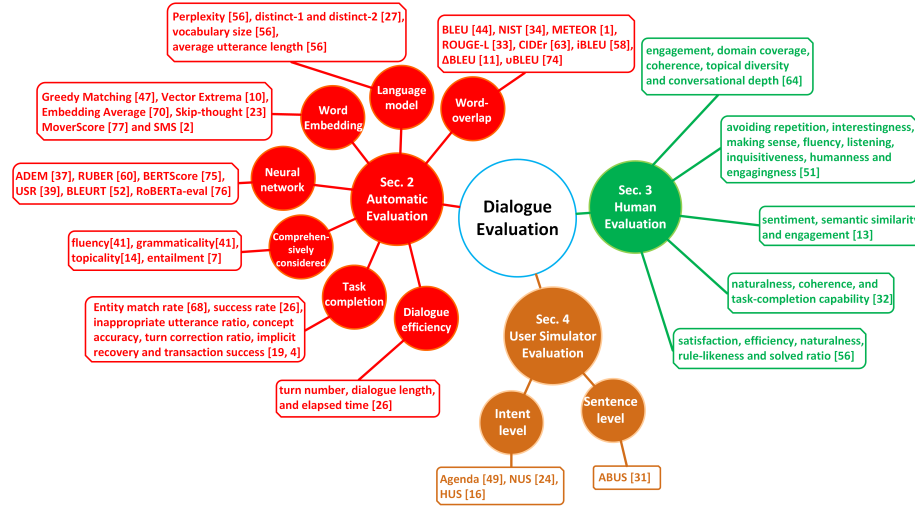


Fig. 2. Summary of evaluation models and metrics. (See Section 2,3,4 for Details)

To balance the cost and accuracy, some hybrid evaluation methods are also proposed to leverage human and automatic evaluation, e.g. Δ BLEU [11], HUSE [17] and GENIE [22].

4 User simulator-based evaluation

Human evaluation is time-consuming and not scalable [55], while automatic evaluation metrics often lead to turn-level and system-level bias. To realize convenient and dynamic evaluation of dialogue systems, user simulators are needed to interact with the dialogue system. To this end, various efforts have been made to build user simulators which mimic human examiners to evaluate dialogue agents through interaction [50].

The commonly adopted method is to generate simulated dialogue by interacting with dialogue system and then assess the reality of the simulated dialogues. Accordingly, the evaluation can be categorized into intent-level and sentence-level. The intent-level focuses on dialogue policy part, interacting with dialogue systems through dialogue act or template-based utterances [49,54,16]. Another line of work evaluates a complete task-oriented dialogue systems with fluent, dynamically generated natural language [31]. López-Cózar et al.[36] tests the performance of spoken dialogue systems by artificially simulating the behavior of three types of user (very cooperative, cooperative and not very cooperative) to interact with previously developed dialogue systems. The user simulator enables the identification of problems relating to the speech recognition, spoken language understanding, and dialogue management components of the system. User simulators can also be used to perform cross validation to measure the miss-distance between automatic metrics with human evaluation [56]. By employing a simulated user in a range of different experimental conditions, sufficient data

can be generated to support a systematic analysis of potential problems and to enable fine-grained tuning of the system [36].

Table 1. A summarization of the benchmarks.

Platform	Feature
PARADISE [67]	comparing dialogue strategies by decoupling task requirements from the dialogue behaviors
PyDial [62]	a dialogue policy specific platform
ParlAI [40]	end-to-end, multiple conversational tasks
ConvLab [25], ConvLab-2 [78]	supporting system-wise simulated evaluation and human evaluation
Plato [43]	single- or multi-party interactions, easy configuration and debug
GENIE [22]	automate and standardize the human evaluation
GEM [12]	various tasks, multilingual, additional metrics allowed

5 Benchmarks

Given that evaluation method often performs differently across tasks and works, the benchmarks can act as testbed to evaluate the latest advances and facilitate the research frontier. We introduce the representative and prevalent benchmarks here, as listed in Table 1.

PARADISE [67] is the first general evaluation framework for spoken dialogue systems. PARADISE supported comparisons among dialogue strategies by decoupling task requirements from the dialogue behaviors. PyDial⁴ implements two success-based evaluator for dialogues evaluation, including a objective success evaluator to compare the constraints and requests the system identifies with the true values, and a subjective task success evaluator to queries the user about the outcome of the dialogue [62]. ParlAI⁵ provides a unified framework with various conversational tasks, including dialogue system, reading comprehension, etc [40]. Convlab⁶ provides a platform for researchers to develop dialogue systems with various architectures or configurations. It supports system-wise simulated evaluation and also provides the interface of AMT platform for human evaluation [25]. ConvLab-2⁷ is the continuator of ConvLab with more powerful architectures and supporting more datasets [78]. On the other hand, Plato⁸ is a abstraction level platform, which is well-designed for easy of understanding and debugging for conversational agents [43]. GENIE⁹ provides a leaderboard for standardizing the human evaluation on text generation systems [22]. GEM¹⁰

⁴ <http://pydial.org>

⁵ <http://parl.ai>

⁶ <http://convlab.github.io>

⁷ <https://github.com/thu-coai/ConvLab-2>

⁸ <https://github.com/uber-research/plato-research-dialogue-system>

⁹ <https://genie.apps.allenai.org>

¹⁰ <https://gem-benchmark.com/>

provides a benchmark integrated with a wide set of tasks and allows the integration of additional metrics to identify the gaps and then prioritize the direction for improvement [12].

6 Open issues and questions

Evaluation is a non-trivial and understudied task. Besides common challenges inherent in the existing evaluation methods, there are some special challenges on evaluation of the dialogue system.

6.1 Reference-free evaluation

Generally, a dialogue system is trained to imitate the ground truth response from dataset. Most evaluation methods compare the generated response with the ground truth. However, there are a diverse of potential responses to a question in a common conversation. [55] points out that multiple ground truth sentences are helpful to improve the correlation between metrics and human evaluation for which provide more adequate and diversified semantic information. However, datasets collecting is also a time-consuming and laborious work. As stated in Section 2.1, the reference-free evaluation method has become a new trend to improve the scalability and generalization ability of evaluation [48,76].

6.2 Interactively dynamic evaluation

Generally, a conversation is a multi-turn interactive process between the user and the dialogue system. In prior evaluation method, the gold dialogue history (standard user utterance and dialogue response from corpus) is fed to dialogue system in each turn without considering the actually generated response before [68,38,72]. Without dynamic interaction, the errors occurred in prior turn would be buried and can't be passed to next turn. That is to say, the turn-level evaluation can't reflect the system-level performance of dialogue system. So dynamic evaluation method is more effective to strengthen the error checking and recovery mechanisms for dialogue system.

6.3 User simulator with fluent natural language

Existing user simulators only focus on dialogue policy part, interacting with dialogue systems through well-structured formal languages [56]. Such user simulators are insufficient to evaluate a complete task-oriented dialogue systems, as they cannot appropriately assess fluent, dynamically generated natural language. With the recent wave of end-to-end dialogue systems [68,26,38,72], the community starts to emphasize on natural utterances and realistic usage of dialogue systems. It would push researchers to build an effective user simulator which realistically evaluate dialogue agents through interacting with them using fluent natural languages.

7 Conclusion

In this paper, we summarize the evaluation methods of dialogue systems and a method-based taxonomy of the existing evaluation metric is attempted here. Even though much work has done to construct an effective evaluation method, it is still a challenging and understudied work to capture all aspects of dialogue response from naturalness and coherence to long-term engagement and flow. Finally, we believe that the redefinition of user simulator is the most promising direction to build an ideal evaluator.

References

1. Banerjee, S., Lavie, A.: METEOR: an automatic metric for MT evaluation with improved correlation with human judgments. In: Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization@ACL, June 29, 2005
2. Clark, E., Celikyilmaz, A., Smith, N.A.: Sentence mover’s similarity: Automatic evaluation for multi-sentence texts. In: 57th ACL (2019)
3. Cummins, R., Rei, M.: Neural multi-task learning in automated assessment. CoRR **1801.06830** (2018)
4. Danieli, M., Gerbino, E.: Metrics for evaluating dialogue strategies in a spoken language system. In: Proceedings of the 1995 AAAI spring symposium on Empirical Methods in Discourse Interpretation and Generation
5. Deriu, J., Rodrigo, Á., Otegi, A., Echegoyen, G., Rosset, S., Agirre, E., Cieliebak, M.: Survey on evaluation methods for dialogue systems. CoRR **1905.04071** (2019)
6. Dusek, O., Novikova, J., Rieser, V.: Referenceless quality estimation for natural language generation. CoRR **1708.01759** (2017)
7. Dziri, N., Kamaloo, E., Mathewson, K.W., Zaiiane, O.R.: Evaluating coherence in dialogue systems using entailment. In: NAACL-HLT, June 2-7, 2019
8. Finch, A.M., Akiba, Y., Sumita, E.: How does automatic machine translation evaluation correlate with human scoring as the number of reference translations increases? In: LREC 2004, May 26-28, 2004
9. Flesch, R.F.: How to write plain English : a book for lawyers and consumers (1979)
10. Forgues, G., Pineau, J., Larchevêque, J.M., Tremblay, R.: Bootstrapping dialog systems with word embeddings. In: Nips, modern machine learning and natural language processing workshop (2014)
11. Galley, M., Brockett, C., Sordoni, A., Ji, Y., Auli, M., Quirk, C., Mitchell, M., Gao, J., Dolan, B.: deltableu: A discriminative metric for generation tasks with intrinsically diverse targets. In: 53rd ACL, July 26-31, 2015
12. Gehrmann, S., Adewumi, T.P., et al.: The gem benchmark: Natural language generation, its evaluation and metrics. arXiv:2102.01672 (2021)
13. Ghandeharioun, A., Shen, J.H., Jaques, N., Ferguson, C., Jones, N., Lapedriza, Á., Picard, R.W.: Approximating interactive human evaluation with self-play for open-domain dialog systems. In: NeurIPS, December 8-14, 2019
14. Guo, F., Metallinou, A., Khatri, C., Raju, A., Venkatesh, A., Ram, A.: Topic-based evaluation for conversational bots. CoRR **1801.03622** (2018)
15. Gupta, P., Mehri, S., Zhao, T., Pavel, A., Eskénazi, M., Bigham, J.P.: Investigating evaluation of open-domain dialogue systems with human generated multiple references. In: SIGdial, September 11-13, 2019

16. Gur, I., Hakkani-Tür, D.Z., Tür, G., Shah, P.: User modeling for task oriented dialogues. 2018 IEEE Spoken Language Technology Workshop (SLT)
17. Hashimoto, T.B., Zhang, H., Liang, P.: Unifying human and statistical evaluation for natural language generation. In: NAACL-HLT, June 2-7, 2019
18. Hastie, H.: Metrics and evaluation of spoken dialogue systems pp. 131–150 (2012)
19. Hirschberg, J., Nakatani, C.H.: A prosodic analysis of discourse segments in direction-giving monologues. In: 34th ACL, 24-27 June 1996
20. Howcroft, D.M., Belz, A., Clinciu, M., Gkatzia, D., Hasan, S.A., Mahamood, S., Mille, S., van Miltenburg, E., Santhanam, S., Rieser, V.: Twenty years of confusion in human evaluation: NLG needs evaluation sheets and standardised definitions. In: 13thINLG, December 15-18, 2020
21. Kannan, A., Vinyals, O.: Adversarial evaluation of dialogue models. CoRR **1701.08198** (2017)
22. Khashabi, D., Stanovsky, G., Bragg, J., Lourie, N., Kasai, J., Choi, Y., Smith, N.A., Weld, D.S.: Genie: A leaderboard for human-in-the-loop evaluation of text generation. arXiv:2101.06561 (2021)
23. Kiros, R., Zhu, Y., Salakhutdinov, R., Zemel, R.S., Urtasun, R., Torralba, A., Fidler, S.: Skip-thought vectors. In: NIPS 2015, December 7-12
24. Kreyssig, F., Casanueva, I., Budzianowski, P., Gasic, M.: Neural user simulation for corpus-based policy optimisation of spoken dialogue systems. In: 19th SIGdial (2018)
25. Lee, S., Zhu, Q., Takanobu, R., Zhang, Z., Zhang, Y., Li, X., Li, J., Peng, B., Li, X., Huang, M., Gao, J.: Convlab: Multi-domain end-to-end dialog system platform. In: 57th ACL, July 28 - August 2, 2019
26. Lei, W., Jin, X., Kan, M.Y., Ren, Z., He, X., Yin, D.: Sequicity: Simplifying task-oriented dialogue systems with single sequence-to-sequence architectures. In: 56th ACL (2018)
27. Li, J., Galley, M., Brockett, C., Gao, J., Dolan, B.: A diversity-promoting objective function for neural conversation models. In: NAACL-HLT, June 12-17, 2016
28. Li, J., Galley, M., Brockett, C., Spithourakis, G.P., Gao, J., Dolan, W.B.: A persona-based neural conversation model. In: 54th ACL, August 7-12, 2016
29. Li, J., Monroe, W., Ritter, A., Jurafsky, D., Galley, M., Gao, J.: Deep reinforcement learning for dialogue generation. In: EMNLP, November 1-4, 2016
30. Li, L., He, H., Williams, J.D.: Temporal supervised learning for inferring a dialog policy from example conversations. In: 2014 IEEE Spoken Language Technology Workshop, December 7-10, 2014
31. Li, X., Chen, Y., Li, L., Gao, J., Celikyilmaz, A.: End-to-end task-completion neural dialogue systems. In: 8th IJCNLP, November 27 - December 1, 2017
32. Li, X., Wang, Y., Sun, S., Panda, S., Liu, J., Gao, J.: Microsoft dialogue challenge: Building end-to-end task-completion dialogue systems. arXiv:1807.11125 (2018)
33. Lin, C.Y.: ROUGE: A package for automatic evaluation of summaries. In: Text Summarization Branches Out (2004)
34. Lin, C., Och, F.J.: Automatic evaluation of machine translation quality using longest common subsequence and skip-bigram statistics. In: 42nd ACL, 2004
35. Liu, C., Lowe, R., Serban, I., Noseworthy, M., Charlin, L., Pineau, J.: How NOT to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation. In: EMNLP, November 1-4, 2016
36. López-Cózar, R., Callejas, Z., McTear, M.F.: Testing the performance of spoken dialogue systems by means of an artificially simulated user. *Artif. Intell. Rev.* **26**(4) (2006)

37. Lowe, R., Noseworthy, M., Serban, I.V., Angelard-Gontier, N., Bengio, Y., Pineau, J.: Towards an automatic turing test: Learning to evaluate dialogue responses. In: 55th ACL (2017)
38. Madotto, A., Wu, C.S., Fung, P.: Mem2seq: Effectively incorporating knowledge bases into end-to-end task-oriented dialog systems. In: 56th ACL
39. Mehri, S., Eskénazi, M.: USR: an unsupervised and reference free evaluation metric for dialog generation. In: 58th ACL, July 5-10, 2020
40. Miller, A.H., Feng, W., Batra, D., Bordes, A., Fisch, A., Lu, J., Parikh, D., Weston, J.: Parlai: A dialog research software platform. In: 2017 EMNLP
41. Novikova, J., Dusek, O., Curry, A.C., Rieser, V.: Why we need new evaluation metrics for NLG. In: EMNLP, September 9-11, 2017
42. Novikova, J., Dusek, O., Rieser, V.: The E2E dataset: New challenges for end-to-end generation. In: 18th SIGdial, August 15-17, 2017
43. Papangelis, A., Namazifar, M., Khatri, C., Wang, Y.C., Molino, P., Tur, G.: Plato dialogue system: A flexible conversational ai research platform. arXiv:2001.06463 (2020)
44. Papineni, K., Roukos, S., Ward, T., Zhu, W.: Bleu: a method for automatic evaluation of machine translation. In: 40th ACL, July 6-12, 2002
45. Peng, B., Li, X., Li, L., Gao, J., Celikyilmaz, A., Lee, S., Wong, K.: Composite task-completion dialogue policy learning via hierarchical deep reinforcement learning. In: EMNLP, September 9-11, 2017
46. Pietquin, O., Hastie, H.F.: A survey on metrics for the evaluation of user simulations. *Knowl. Eng. Rev.* **28**(1) (2013)
47. Rus, V., Lintean, M.C.: A comparison of greedy and optimal assessment of natural language student input using word-to-word similarity metrics. In: Proceedings of the 7th Workshop on Building Educational Applications Using NLP, June 7, 2012
48. Sai, A.B., Gupta, M.D., Khapra, M.M., Srinivasan, M.: Re-evaluating ADEM: A deeper look at scoring dialogue responses. In: AAAI, Jan 27 - Feb 1, 2019
49. Schatzmann, J., Thomson, B., Weilhammer, K., Ye, H., Young, S.: Agenda-based user simulation for bootstrapping a pomdp dialogue system. In: NAACL-HLT 2007
50. Schatzmann, J., Weilhammer, K., Stuttle, M., Young, S.: A survey of statistical user simulation techniques for reinforcement-learning of dialogue management strategies. *The knowledge engineering review* **21**(2) (2006)
51. See, A., Roller, S., Kiela, D., Weston, J.: What makes a good conversation? how controllable attributes affect human judgments. In: NAACL-HLT, June 2-7, 2019
52. Sellam, T., Das, D., Parikh, A.P.: Bleurt: Learning robust metrics for text generation. In: 58th ACL (2020)
53. Serban, I.V., Sordoni, A., Lowe, R., Charlin, L., Pineau, J., Courville, A.C., Bengio, Y.: A hierarchical latent variable encoder-decoder model for generating dialogues. In: 31st AAAI, February 4-9, 2017
54. Shah, P., Hakkani-Tür, D., Liu, B., Tür, G.: Bootstrapping a neural conversational agent with dialogue self-play, crowdsourcing and on-line reinforcement learning. In: NAACL-HLT, June 1-6, 2018
55. Sharma, S., Asri, L.E., Schulz, H., Zumer, J.: Relevance of unsupervised metrics in task-oriented dialogue for evaluating natural language generation. arXiv:1706.09799 (2017)
56. Shi, W., Qian, K., Wang, X., Yu, Z.: How to build user simulators to train rl-based dialog systems. In: 2019 EMNLP-IJCNLP
57. Su, P., Gasic, M., et al.: On-line active reward learning for policy optimisation in spoken dialogue systems. In: 54th ACL, August 7-12, 2016

58. Sun, H., Zhou, M.: Joint learning of a dual SMT system for paraphrase generation. In: 50th ACL, July 8-14, 2012
59. Takanobu, R., Zhu, Q., Li, J., Peng, B., Gao, J., Huang, M.: Is your goal-oriented dialog model performing really well? empirical analysis of system-wise evaluation. In: Proceedings of the 21st SIGDial (2020)
60. Tao, C., Mou, L., Zhao, D., Yan, R.: RUBER: an unsupervised method for automatic evaluation of open-domain dialog systems. In: 32nd AAAI, Feb 2-7, 2018
61. Ultes, S., Budzianowski, P., et al.: Domain-independent user satisfaction reward estimation for dialogue policy learning. In: Interspeech 2017, August 20-24, 2017
62. Ultes, S., Rojas-Barahona, L.M., et al.: Pydial: A multi-domain statistical dialogue system toolkit. In: ACL 2017, July 30 - August 4
63. Vedantam, R., Zitnick, C.L., Parikh, D.: Cider: Consensus-based image description evaluation. In: CVPR, June 7-12, 2015
64. Venkatesh, A., Khatri, C., Ram, A., Guo, F., Gabriel, R., Nagar, A., Prasad, R., Cheng, M., Hedayatnia, B., Metallinou, A., et al.: On evaluating and comparing open domain dialog systems. arXiv:1801.03625 (2018)
65. Vinyals, O., Le, Q.V.: A neural conversational model. CoRR **1506.05869** (2015)
66. Walker, M.A., Kamm, C.A., Litman, D.J.: Towards developing general models of usability with PARADISE. Nat. Lang. Eng. **6**(3&4) (2000)
67. Walker, M.A., Litman, D.J., Kamm, C.A., Abella, A.: PARADISE: A framework for evaluating spoken dialogue agents. In: 35th ACL-EACL
68. Wen, T., Vandyke, D., Mrkšić, N., Gašić, M., Rojas-Barahona, L., Su, P., Ultes, S., Young, S.: A network-based end-to-end trainable task-oriented dialogue system. In: EACL 2017
69. Wen, T., Gasic, M., Mrksic, N., Su, P., Vandyke, D., Young, S.J.: Semantically conditioned lstm-based natural language generation for spoken dialogue systems. In: EMNLP, September 17-21, 2015
70. Wieting, J., Bansal, M., Gimpel, K., Livescu, K.: Towards universal paraphrastic sentence embeddings. In: 4th ICLR, May 2-4, 2016
71. Williams, J.D., Asadi, K., Zweig, G.: Hybrid code networks: practical and efficient end-to-end dialog control with supervised and reinforcement learning. In: 55th ACL, July 30 - August 4, 2017
72. Wu, C.S., Socher, R., Xiong, C.: Global-to-local memory pointer networks for task-oriented dialogue. In: 7th ICLR (2019)
73. Young, S., Gašić, M., Thomson, B., Williams, J.D.: Pomdp-based statistical spoken dialog systems: A review. Proceedings of the IEEE **101**(5) (2013)
74. Yuma, T., Yoshinaga, N., Toyoda, M.: ubleu: Uncertainty-aware automatic evaluation method for open-domain dialogue systems. In: 58th ACL: SRW (2020)
75. Zhang, T., Kishore, V., Wu, F., Weinberger, K.Q., Artzi, Y.: Bertscore: Evaluating text generation with BERT. In: ICLR, April 26-30, 2020
76. Zhao, T., Lala, D., Kawahara, T.: Designing precise and robust dialogue response evaluators. In: 58th ACL, July 5-10, 2020
77. Zhao, W., Peyrard, M., Liu, F., Gao, Y., Meyer, C.M., Eger, S.: Moverscore: Text generation evaluating with contextualized embeddings and earth mover distance. In: EMNLP-IJCNLP 2019
78. Zhu, Q., Zhang, Z., Fang, Y., Li, X., Takanobu, R., Li, J., Peng, B., Gao, J., Zhu, X., Huang, M.: Convlab-2: An open-source toolkit for building, evaluating, and diagnosing dialogue systems. In: 58th ACL, July 5-10, 2020